

---

# HUMAN-CENTERED METRICS FOR DIALOG SYSTEM EVALUATION

---

**Salvatore Giorgi**  
University of Pennsylvania  
sgiorgi@sas.upenn.edu

**Shreya Havaldar**  
University of Pennsylvania  
shreyah@seas.upenn.edu

**Farhan Ahmed**  
Stony Brook University  
farhaahmed@cs.stonybrook.edu

**Zuhaib Akhtar**  
New York University  
za2023@nyu.edu

**Shalaka Vaidya**  
New York University  
sav5821@nyu.edu

**Gary Pan**  
New York University  
hp2190@nyu.edu

**Lyle H. Ungar**  
University of Pennsylvania  
ungar@cis.upenn.edu

**H. Andrew Schwartz**  
Stony Brook University  
has@cs.stonybrook.edu

**João Sedoc**  
New York University  
jsedoc@stern.nyu.edu

## ABSTRACT

We present metrics for evaluating dialog systems through a psychologically-grounded “human” lens: conversational agents express a diversity of both *states* (short-term factors like emotions) and *traits* (longer-term factors like personality) just as people do. These interpretable metrics consist of five measures from established psychology constructs that can be applied both across dialogs and on turns within dialogs: emotional entropy, linguistic style and emotion matching, as well as agreeableness and empathy. We compare these human metrics against 6 state-of-the-art automatic metrics (e.g. BARTScore and BLEURT) on 7 standard dialog system data sets. We also introduce a novel data set, the Three Bot Dialog Evaluation Corpus, which consists of annotated conversations from ChatGPT, GPT-3, and BlenderBot. We demonstrate the proposed human metrics offer novel information, are uncorrelated with automatic metrics, and lead to increased accuracy beyond existing automatic metrics for predicting crowd-sourced dialog judgements. The interpretability and unique signal of our proposed human-centered framework make it a valuable tool for evaluating and improving dialog systems.

## 1 Introduction

Open-domain dialog systems are typically evaluated with automatic lexico-semantic metrics and human judgements, which have a number of drawbacks [1]. Automatic methods (e.g. BLEU, METEOR, and ROUGE), often rely heavily on overlap and fail to capture the diversity of dialog systems [2]. This typically results in small associations with human judgements [2, 3]. On the other hand, human judgements are expensive to scale and lack standardization [4, 5, 6]. Automatic metrics that capture human-like conversation agent attributes could drive dialog system improvements.

In this work, we propose two classes of psychologically-grounded metrics for evaluating open-domain dialog systems as if they were human, taking queues from Giorgi et al. [7] which characterized Twitter spambots through a number of human states and traits. We propose three general classes of psychologically-grounded measures: (1) *states* (changing within a dialog, such as emotion), (2) *traits* (slower to change, such as personality), and (3) linguistic matching (i.e., how well chatbots match the linguistic cues of the other entity in the conversation). Along with these measures, we propose a hierarchical framework for evaluating dialog systems, where turns happen with dialogues, dialogues are nested with agents, and agents are nested with the larger dialog system. We also introduce a benchmark dialog data set of conversations with ChatGPT, GPT-3, and BlenderBot, annotated at both the turn- and dialog-level. Finally, we

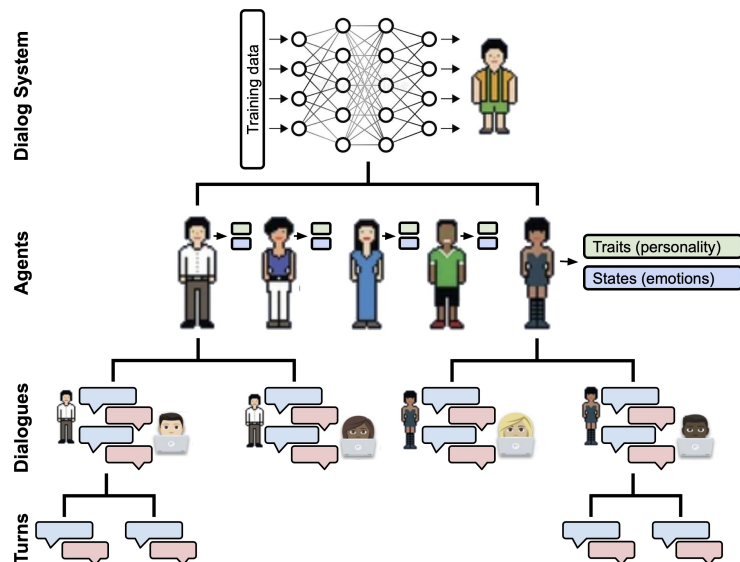


Figure 1: The proposed hierarchical structure of open-domain dialog systems. Here Turns (utterances) are nested within Dialogues (a complete exchange), which are nested within Agents (trained instances of a dialog system), which are nested within Dialog Systems (the overall system architecture). We note that Turns consist of one utterance from both the agent and entity, but turn level outcomes do not need to be measured with the agent’s utterance preceding the entity’s (as depicted in the figure).

systematically compare these human-centered metrics against an extensive list of automatic metrics (6) on the above mentioned data set as well as seven additional publicly available data sets.

**Contributions** Our contributions include: (1) three classes of psychologically grounded metrics with five specific instances of metrics within these three classes, (2) a hierarchical framework for dialog system evaluations, (3) a publicly available data set of conversations from state-of-the-art dialog systems (ChatGPT, GPT-3, and BlenderBot) with turn- and dialog-level annotations, and (4) a systematic evaluation against 6 existing metrics across 7 data sets. We show that (a) the psychological metrics are uncorrelated with the automatic methods and (b) the psychological metrics increase accuracy when combined with the automatic methods to predict crowd-sourced dialog system judgements. Together, these results show that our psychological metrics can be used in tandem with existing metrics to further characterize dialog systems.

## 2 Related Work

There is a growing set of methods to embed language processing within human contexts [8, 9, 10]. Most of such work has focused on the modeling side rather than evaluations, for example, *creating* agents with human-like traits such as empathy [11], trust [12], emotion [13, 14], and personalizations and personas [15, 16, 17, 18]. On the other hand, few have attempted to *evaluate* dialog agents as human with a number of human-like metrics. Adiwardana et al. [19] proposed a metric which jointly measures “making sense” and being specific, both basic and important attributes of conversations. More directly, some have quantified “humanness” subjectively through crowd-sourcing: “Which speaker sounds more human?” [20, 18, 21].

A parallel line of work seeks to improve language models by making them more human-aligned. Santurkar et al. [22] evaluates whose opinions language models reflect via public opinion polls and Binz et al. [23] assesses whether language models reflect the cognitive ability of humans. Glaese et al. [24] establishes rules to make dialog agents more helpful and harmless. Additional work on assessing alignment of agents [25, 26, 27] focuses on measuring and minimizing the attributes of agents that make them bad conversationalists (hate speech, toxicity, controversy, etc).

Our work takes a different step toward evaluation of “human-like” dialog. We propose three classes of psychologically grounded and human-centered measures which can be used to evaluate dialog systems. These metrics additionally seek to measure and exemplify the attributes of dialog agents that make them *good* conversationalists. We see this as a step toward answering the call for a human-like open-domain system [19], and for integrating current steps toward

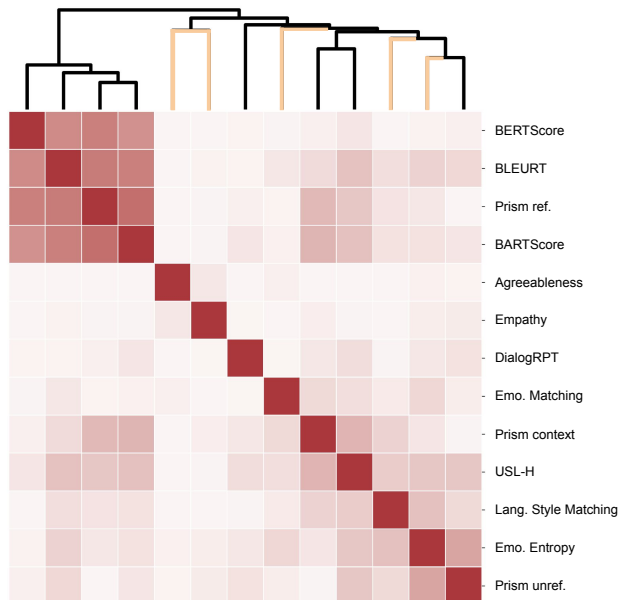


Figure 2: Correlations between psychological and automatic metrics. We cluster both the rows and columns based on absolute correlations. Psychological metrics denoted with orange dendrogram lines.

this.<sup>1</sup> Despite the general structure of our proposed metrics, we note that a number of dialog systems are task or goal oriented, such as question/answer systems [28] or systems designed for highly specific tasks such as trip planning [29] and customer service [30]. Such systems may be considered outside of the scope of our formulation, in that scheduling a trip is fundamentally different from, for example, a conversational chatbot related to COVID-19 vaccines, which may need additional social and cultural context.

### 3 The Human-Centered Hierarchical Framework for Dialogue Systems

We propose that dialog systems should be evaluated across all four levels with considerations for both models and open source data sets (e.g., there should be sufficient information to group turns within dialogues and dialogues within agents). Figure 1 shows our proposed four part hierarchical framework. Each part is defined as follows:

**Dialog System** This is the overall architecture of the system and the top level of our hierarchy. The dialog system produces agents and could have the ability to produce agents across a spectrum of desirable attributes (i.e., tunable).

**Agent** An agent is a specific instance of a dialog system and we note that a single dialog system can produce a number of agents. For example, agents could be trained instances of a specific dialog system architecture, where each instance uses different training data. Other examples could include instances of dialog systems with tunable parameters (e.g., empathy) where each agent has a different instantiated parameter value (e.g., one agent with low empathy, another with high empathy). With this view, a dialog system can be thought of as an agent generator.

**Dialog** A dialog is a complete back and forth exchange between an agent and another entity (another dialog agent or a human). A single agent can engage in multiple dialogues.

**Turn** A turn is a specific utterance with a dialog. This could include the second entity’s preceding or proceeding response.

### 4 Classes of Human-like Measures

We propose two classes of measures: (1) states and traits and (2) linguistic matching, rooted in fundamental psychological measurements of humans and their social relationships and interactions (i.e., linguistic matching). The next section operationalizes these classes across five metrics.

<sup>1</sup>For example, Roller et al. [18] propose evaluating both “engaging talking points” and “consistent persona” which are captured within our proposed metrics via state and trait metrics, respectively, where consistency can be measured across multiple dialogues.

**States and Traits** The state vs. trait distinction is ubiquitous in psychology, with a long history [31]. A standard textbook defines *state measures* as thoughts, feelings, and behaviors in a specific place and time. *Trait measures* are those which generalize across situations that are stable over time, and systematically differ across people [32]. Emotions are states while personalities are traits. In relation to standard NLP tasks, past work has found stance-detection being more trait-like while sentiment is a more state-like outcome [10]. It is important to distinguish the measures we use (e.g., personality), grounded against validated psychological instruments, from proxies for these constructs used in other works (e.g., personas). While proxy measures such as “likes” correlate with personality [33], they are not direct assessments of the constructs.

Within our hierarchy, we propose that states and traits be considered as follows: First, at the top level, dialog systems should have the capacity to produce a number of agents with varying traits, with each agent maintaining its given traits across dialogues. One should thus be able to measure variation in traits across agents from a single dialog system and stability in each agent across multiple dialogues. On the other hand, states should vary within dialogues and agents should have the capacity to exhibit a range of states.

**Linguistic Matching** Linguistic matching has been observed in many settings, and has been shown to predict power differentials [34], relationship stability [35], cooperation [36], and empathy ratings of therapists [37]. More generally, the psycholinguistic theory of communication accommodation has studied such unconscious matching tendencies in postures, facial expressions, pitch, pausing, length, and use of function words [38]. Besides sentence embedding similarity [39], to our knowledge, such extensive matching phenomena have yet to be fully studied in open-domain dialog systems, despite being applied in other NLP settings [40, 41].

As measured within our proposed hierarchy, linguistic matching is a property of dialogues and turns. For example, one could measure function word matching in a single turn (how well does the agent match the prompt?) or across a dialog.

## 5 Psychological Metrics

*Psychological metrics* operationalize the human-like measures using models trained on other data sets to predict e.g. emotion and personality. Because they were not specifically designed for evaluating dialog systems, they are not optimized to correlate with the gold standard human judgements in the data sets (e.g., appropriateness). Five metric scores were produced at the turn and dialog level (depending on the data set) and then correlated with a number of crowd-sourced evaluations:<sup>2</sup>

**Emotional Entropy** Using the NRC Hashtag Emotion Lexicon [42] we estimate Plutchik’s eight basic emotions: anger, anticipation, disgust, fear, joy, sadness, surprise, and trust [43]. This emotion lexicon, which is a set of weighted words for each emotion category, was automatically derived over tweets with emotion hashtags (e.g., *#anger* and *#joy*). The lexicon is applied to every observation in each data set (i.e., we summed weighted word frequencies which were weighted according to their weight within each emotion category) and then the entropy of the normalized emotion vector is calculated. Emotions (and, thus, emotional entropy) are state measures and can be estimated at multiple levels of the hierarchy: turn, dialog, and agent.

**Agreeableness** We used a language based personality model to estimate the agreeableness dimension of the Big Five personality traits [44]. This model had an out-of-sample prediction accuracy (Pearson  $r$ ) of .35 and was built over 1-3grams and 2,000 LDA topics (Latent Dirichlet Allocation; [45]). Thus, for each dialog, we extracted 1-3grams and loadings for the 2,000 LDA topics and applied the pre-trained regression model, which produced an agreeableness score for each observation. We include agreeableness in our final five metrics since it out performed the other four personality measures (openness to experience, conscientiousness, extraversion, and neuroticism) on the test data. Agreeableness (and personality, in general) is a trait measure that would typically be defined at the agent level or above (e.g., for a given dialog system, does agreeableness vary across agents and is it stable within an agent), though do to lack of agent-level data in the task we estimate agreeableness for dialog-level data sets.

**Empathy** We build a model to predict empathy, as measured by the Empathic Concern subscale of the Interpersonal Reactivity Index (IRI) [46]. We use an existing empathy data set [47, 48] and build a model over 2,805 participants who shared their Facebook status data and answered the IRI questionnaire. Using 10-fold cross validation, we predicted the empathic concern scores from a Ridge penalized linear regression using the same set of 2,000 LDA topics described above. The final model resulted in an out-of-sample Pearson  $r$  of 0.26. In order to obtain Empathic Concern estimates for each dialog, we extracted 2,000 LDA topic loadings for each observation and applied the pre-trained regression

---

<sup>2</sup>See Appendix for full details on automatic metrics, data sets, and crowd-source annotations.

model. Empathic Concern is a trait level measure. Similar to agreeableness, this would typically be defined at the agent level or above, but for this task we estimate Empathic Concern for dialogues.

**Language Style Matching** We use the definition provided by Ireland et al. [35]: 1 minus the normalized absolute difference in function word use between the agent and entity. This score was calculated for nine separate function word categories in the Linguistic Inquiry and Word Count (LIWC) dictionary [49]: personal pronouns, impersonal pronouns, articles, conjunctions, prepositions, auxiliary verbs, high frequency adverbs, negations, and quantifiers. Turn and dialog level scores were averaged across the nine categories. This is a form of Linguistic Matching which can be measured at the turn, dialog, and agent levels.

**Emotion Matching** Again, we use the NRC Hashtag Emotion Lexicon [42] and calculate the Spearman rank correlation between the agents emotions and the prompts. Emotion Matching is a form of Linguistic Matching which can be measured at the turn, dialog, and agent levels.

## 6 Data

To evaluate our human metrics we collect a novel data set, the Three Bot Dialog Evaluation Corpus, from three state-of-the-art dialog systems and evaluate the dialogues at both the turn and dialog level via a crowd sourcing (Amazon Mechanical Turk). We also evaluate our metrics on several additional open-source data sets, the DSTC10 Track 5 Test Corpus.

**Three Bot Dialog Evaluation Corpus** Here we introduce the Three Bot Dialog Evaluation Corpus (TBD-Q1-2023 OR TBD; Quarter 1 of 2023). This data set consists of conversations with three chatbots: ChatGPT, GPT-3 [50], and BlenderBot [51]. For each chatbot, we collected 21 dialogues with an average of 14.6 turns per dialogue.

We then collect human judgements at both the turn and dialog level for each conversation in the data set. At the turn level, we ask crowd workers to evaluate the appropriateness, content, grammar, and relevance. At the dialog level, we ask crowd workers to evaluate the conversation at for coherence, informativeness, likability, and overall. The exact evaluation question text is included in the Appendix.

**DSTC10 Track 5 Test Corpus** In order to further evaluate our human metrics, we use a test corpus from The Tenth Dialog System Technology Challenge (DSTC10) Track 5 Automatic Evaluation and Moderation of Open-domain Dialogue Systems [1]. This evaluation data set combined 7 *turn level* data sets into a single data set: DailyDialog, PersonaChat, Topical-DTSC10, Persona-DSTC10, TopicalChat-USR, PersonaChat-USR, and DailyDialog.<sup>2</sup> Since this data set is available at the turn level, we evaluate our three turn level metrics: emotional entropy, emotion matching, and language style matching.

## 7 Evaluation

**Automatic Metrics** We use 6 common dialog system metrics: BARTScore [52], BERTScore [53], BLEURT [54], DialogRPT [55], Prism [56], Mauve [57], and USL-H [58].<sup>2</sup> TBD-Q1-2023 is evaluated using DialogRPT, Mauve, and USL-H only. All metrics, with the exception of Mauve, are used to evaluate the DSTC10 Track 5 Test Corpus. As these are all turn-level metrics, we average metrics across turns to create dialog-level scores when evaluating TBD-Q1-2023 at the dialog-level.

**Evaluation** We create three models which contain varying sets of independent variables: (1) the automatic metric (“A”), (2) the psychological metric (“P”), and (3) both the psychological and automatic metrics together (“P+A”). In all models, the dependent variable is the average crowd-source annotation.<sup>2</sup> Additionally, all variables are mean centered and standardized, so that the resulting standard deviation is equal to 1. We report model fit via adjusted  $R^2$ . We also perform a paired t-test between the mean absolute residuals of the “A” and “P+A” models to see if the psychological metrics add significant predictive value above the automatic metrics alone. We then apply a Benjamini-Hochberg False Discovery Rate correction to compensate for the large number of comparisons [59].

## 8 Results

Figure 2 shows the clustered correlations between the *psychological* and *automatic* metrics on the DSTC10 Track 5 Test Corpus. Three distinct clusters appear: (1) BARTScore, BERTScore, BLEURT, and Prism ref.; (2) empathy and

agreeableness, which are the two proposed non-turn level metrics (though, here, they are applied at the turn level); and (3) the remaining metrics. As expected, all reference based contextualized embedding methods should cluster together.

		Automatic Metric Alone	Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
			P	P+A	P	P+A	P	P+A	P	P+A
TBD	DialogRPT	.133	.017	.138	.014	.138	.001	.135	.031	.144**
	Mauve	-.001	.017	.016*	.014	.013*	.001	.000	.031	.031**
	USL-H	.000	.017	.017*	.014	.014**	.001	.001	.031	.031***
DSTC10	BARTScore	.085	.008	.086**	.008	.089***	.004	.085	.009	.090***
	BERTScore	.062	.008	.066***	.008	.068***	.004	.064*	.009	.072***
	BLEURT	.081	.008	.082	.008	.086*	.004	.082	.009	.087*
	DialogRPT	.006	.008	.014***	.008	.014***	.005	.011***	.009	.021***
	Prism ref.	.051	.008	.054***	.008	.057***	.004	.052	.009	.059***
	USL-H	.105	.008	.105	.008	.107**	.004	.105	.009	.106**

Table 1: Turn-level Results: Reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metrics as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

Table 1 shows the comparison between the psychological and automatic metrics when predicting the turn level human judgements. A number of the state-of-the-art automatic metrics performed well, such as BARTScore and BLEURT. While the psychological metrics did not perform as well, we see that emotional entropy, emotion matching, and language style matching all increase predictive accuracy when combined with the automatic metrics. Table 2 shows the results of the dialogue-level analysis, predicting the Overall annotation. (See Appendix Tables 3, 4, and 5 for coherence, informativeness, and likability results.)

		Automatic Metric Alone	Agreeableness		Empathy		All Psych.	
			P	P+A	P	P+A	P	P+A
TBD	DialogRPT	.180	.094	.175	.031	.168	.089	.161
	Mauve	.091	.094	.120	.031	.084	.089	.106
	USL-H	.002	.094	.120	.031	.089	.089	.146

Table 2: Dialogue-level results predicting the Overall rating: Reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metrics as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables.

Taken together, the psychological metrics were not highly predictive alone when compared to state-of-the-art metrics (which is expected since the psychological metrics are not specific for dialog evaluations), yet they are capturing unique relevant signal for dialog quality. Similar results hold across an additional 10 out of 12 open-domain dialog evaluation data sets in the Appendix.<sup>2</sup>

## 9 Conclusions

In this paper, we proposed a hierarchical framework for evaluating open-domain dialog systems with human-centered measures which consider both trait and state trade-offs (standard measures of human constructs) and linguistic matching (indicators of social relationships and interactions). Five metrics were evaluated, which examined trait level features (agreeableness and empathy), state level variation (emotional entropy), and linguistic matching (style and emotion matching), and compared to state-of-the-art automatic metrics. Across multiple data sets we show that the psychological metrics (1) do not correlate with automatic metrics and (2) increase accuracy when predicting gold standard human judgements, showing that the psychological metrics are picking up on unique signal when evaluating open-domain dialog systems.

## Ethical Considerations

There are a number of ethical considerations when constructing and evaluation dialog systems, many of which have been outlined by Roller et al. [18]. These include privacy (since online dialog may contain sensitive information), toxic

and offensive content, and, on the part of the researcher, openness to sharing findings. With regard to the current work, imparting system with human qualities such as personality and socio-demographics must be handled with the utmost sensitivity. Biases in training data, misclassifications in downstream tasks, and reliance on outdated social constructs (i.e., binary gender) are just a few examples of how automated systems can fail and further marginalize vulnerable populations [60, 61, 62]. Specifically, the models used in this study (e.g., empathy and agreeableness) are trained on majority U.S. and monolingual English speaking populations and may fail to generalize to minority or non-US populations. On the other hand, the alternative also suffers from similar concerns, namely that dialog systems may exhibit extremely limited variation in such traits. One could imagine a similar situation to the so-called “Wall Street Journal effect” (i.e., part-of-speech taggers are only accurate when applied to language written by white men; [63]), where dialog system only converse like middle aged white men. Within our proposed framework, dialog systems should produce agents along a spectrum of such trait level constructs.

It is also important to note that while the proposed classes of metrics (i.e., states/traits and linguistic matching) may be desirable in the context of “human-like” measures, the examples used in the paper (e.g., agreeableness) may not. When presented with a toxic prompt, an agreeable or style matching dialog system will only reinforce the toxicity by agreeing with or matching the prompt, while embedding systems with social norms may help alleviate such issues [64]. In general, more human-like dialog systems, as enabled by this approach, can be used both for good (better support for mental health) and for evil (more effective deception and misinformation). Thus, care must be taken when choosing constructs to be embedded in dialog systems.

### Limitations

While we have attempted to evaluate our metrics on a large number of public data sets and compare them against many state-of-the-art metrics, there are a number of limitations. First, our proposed hierarchy is defined at four levels, yet most data sets used in this paper contain human judgements at the turn level only. To our knowledge, no public data sets contain human judgements at all of the level proposed in the hierarchy. Second, the psychological metrics are not developed for dialog system evaluations and may fail to capture the nuances of this domain. For example, the agreeableness model was trained lifetime post histories from Facebook users, and thus one may not expect this to work well on short responses within a dialog or even conversations in general. Next, the specific metrics proposed in this paper (e.g., agreeableness and empathy) are just five examples of psychologically grounded measures which could be applied in this setting. We do not claim to have attempted a thorough investigation across all possible (or even a large number of) psychological metrics. Finally, there is no reason to expect the proposed psychological metrics to correlate with the human judgements. For example, it is not immediately clear that emotional entropy should correlate with either “appropriateness” and “relevance”.

### References

- [1] Zhang Chen, João Sedoc, Luis Fernando D’Haro, Rafael Banchs, and Alexander Rudnicky. Automatic evaluation and moderation of open-domain dialogue systems. 2021.
- [2] Chia-Wei Liu, Ryan Lowe, Iulian Vlad Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132, 2016.
- [3] Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echevoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review*, 54(1):755–810, 2021.
- [4] João Sedoc, Daphne Ippolito, Arun Kirubarajan, Jai Thirani, Lyle Ungar, and Chris Callison-Burch. Chateval: A tool for chatbot evaluation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 60–65. Association for Computational Linguistics, 2019.
- [5] David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland, December 2020. Association for Computational Linguistics.
- [6] Eric Smith, Orion Hsu, Rebecca Qian, Stephen Roller, Y-Lan Boureau, and Jason Weston. Human evaluation of conversations is an open problem: comparing the sensitivity of various methods for evaluating dialogue agents.

- In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 77–97, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [7] Salvatore Giorgi, Lyle Ungar, and H. Andrew Schwartz. Characterizing social spambots by their human traits. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5148–5158, Online, August 2021. Association for Computational Linguistics.
- [8] Svitlana Volkova, Theresa Wilson, and David Yarowsky. Exploring demographic language variations to improve multilingual sentiment analysis in social media. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1815–1827, 2013.
- [9] Dirk Hovy. Demographic factors improve classification performance. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 752–762, 2015.
- [10] Veronica Lynn, Salvatore Giorgi, Niranjan Balasubramanian, and H. Andrew Schwartz. Tweet classification without the tweet: An empirical examination of user versus document attributes. In *Proceedings of the Third Workshop on Natural Language Processing and Computational Social Science*, pages 18–28, Minneapolis, Minnesota, 2019. Association for Computational Linguistics.
- [11] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, 2019.
- [12] David Novick, Mahdokht Afravi, Adriana Camacho, Laura J Hinojos, and Aaron E Rodriguez. Inducing rapport-building behaviors in interaction with an embodied conversational agent. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, pages 345–346. ACM, 2018.
- [13] Xianda Zhou and William Yang Wang. MojiTalk: Generating emotional responses at scale. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1128–1137, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [14] Bernd Huber, Daniel McDuff, Chris Brockett, Michel Galley, and Bill Dolan. Emotional dialogue generation using image-grounded language models. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 277. ACM, 2018.
- [15] Jiwei Li, Michel Galley, Chris Brockett, Georgios P Spithourakis, Jianfeng Gao, and Bill Dolan. A persona-based neural conversation model. *arXiv preprint arXiv:1603.06155*, 2016.
- [16] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*, 2018.
- [17] Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. Training millions of personalized dialogue agents. *arXiv preprint arXiv:1809.01984*, 2018.
- [18] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, et al. Recipes for building an open-domain chatbot. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, 2021.
- [19] Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*, 2020.
- [20] Margaret Li, Jason Weston, and Stephen Roller. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*, 2019.
- [21] Jan Deriu, Don Tuggener, Pius von Däniken, Jon Ander Campos, Alvaro Rodrigo, Thiziri Belkacem, Aitor Soroa, Eneko Agirre, and Mark Cieliebak. Spot the bot: A robust and efficient framework for the evaluation of conversational dialogue systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3971–3984, Online, November 2020. Association for Computational Linguistics.
- [22] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cino Lee, Percy Liang, and Tatsunori Hashimoto. Whose opinions do language models reflect?, 2023.
- [23] Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. *Proceedings of the National Academy of Sciences*, 120(6), feb 2023.
- [24] Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz,



- Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. Improving alignment of dialogue agents via targeted human judgements, 2022.
- [25] Amanda Askill, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for alignment, 2021.
- [26] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askill, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- [27] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askill, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022.
- [28] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, 2017.
- [29] Layla El Asri, Hannes Schulz, Shikhar Kr Sarma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. Frames: a corpus for adding memory to goal-oriented dialogue systems. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 207–219, 2017.
- [30] Lei Cui, Shaohan Huang, Furu Wei, Chuanqi Tan, Chaoqun Duan, and Ming Zhou. Superagent: A customer service chatbot for e-commerce websites. In *Proceedings of ACL 2017, System Demonstrations*, pages 97–102, 2017.
- [31] HA Carr and FA Kingsbury. The concept of traits. *Psychological Review*, 45(6):497, 1938.
- [32] Virgil Zeigler-Hill and T Shackelford. Encyclopedia of personality and individual differences. 2020.
- [33] Michal Kosinski, David Stillwell, and Thore Graepel. Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the national academy of sciences*, 110(15):5802–5805, 2013.
- [34] Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang, and Jon Kleinberg. Echoes of power: Language effects and power differences in social interaction. In *Proceedings of the 21st international conference on World Wide Web*, pages 699–708, 2012.
- [35] Molly E. Ireland, Richard B. Slatcher, Paul W. Eastwick, Lauren E. Scissors, Eli J. Finkel, and James W. Pennebaker. Language style matching predicts relationship initiation and stability. *Psychological Science*, 22(1):39–44, 2011. PMID: 21149854.
- [36] Joseph H Manson, Gregory A Bryant, Matthew M Gervais, and Michelle A Kline. Convergence of speech rate in conversation predicts cooperation. *Evolution and Human Behavior*, 34(6):419–426, 2013.
- [37] Sarah Peregrine Lord, Elisa Sheng, Zac E Imel, John Baer, and David C Atkins. More than reflections: Empathy in motivational interviewing includes language style synchrony between therapist and client. *Behavior therapy*, 46(3):296–303, 2015.
- [38] Howard Ed Giles, Justine Ed Coupland, and Nikolas Ed Coupland. Contexts of accommodation: Developments in applied sociolinguistics. 1991.
- [39] Chen Zhang, Luis Fernando D’Haro, Rafael E Banchs, Thomas Friedrichs, and Haizhou Li. Deep am-fm: Toolkit for automatic dialogue evaluation. In *Conversational Dialogue Systems for the Next Decade*, pages 53–69. Springer, 2021.
- [40] Cristian Danescu-Niculescu-Mizil, Michael Gamon, and Susan Dumais. Mark my words! linguistic style accommodation in social media. In *Proceedings of the 20th international conference on World wide web*, pages 745–754, 2011.
- [41] Cristian Danescu-Niculescu-Mizil and Lillian Lee. Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs. In *Proceedings of the 2nd Workshop on Cognitive*

- Modeling and Computational Linguistics*, pages 76–87, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.
- [42] Saif M Mohammad and Svetlana Kiritchenko. Using hashtags to capture fine emotion categories from tweets. *Computational Intelligence*, 31(2):301–326, 2015.
- [43] Robert Plutchik. A general psychoevolutionary theory of emotion. In *Theories of emotion*, pages 3–33. Elsevier, 1980.
- [44] Gregory Park, H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Michal Kosinski, David J Stillwell, Lyle H Ungar, and Martin EP Seligman. Automatic personality assessment through social media language. *Journal of personality and social psychology*, 108(6):934, 2015.
- [45] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- [46] Mark H Davis. Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of personality and social psychology*, 44(1):113, 1983.
- [47] Muhammad Abdul-Mageed, Anneke Buffone, Hao Peng, Salvatore Giorgi, Johannes C Eichstaedt, and Lyle H Ungar. Recognizing pathogenic empathy in social media. In *ICWSM*, pages 448–451, 2017.
- [48] David B. Yaden, Salvatore Giorgi, Matthew Jordan, Anneke Buffone, Johannes C. Eichstaedt, H. Andrew Schwartz, Lyle H. Ungar, and Paul Bloom. Characterizing empathy and compassion using computational linguistic analysis. *Emotion*, 2023.
- [49] James W Pennebaker, Martha E Francis, and Roger J Booth. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, 71(2001):2001, 2001.
- [50] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [51] Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, et al. Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188*, 2022.
- [52] Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277, 2021.
- [53] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- [54] Thibault Sellam, Dipanjan Das, and Ankur Parikh. Bleurt: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, 2020.
- [55] Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, and William B Dolan. Dialogue response ranking training with large-scale human feedback data. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 386–395, 2020.
- [56] Brian Thompson and Matt Post. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online, November 2020. Association for Computational Linguistics.
- [57] Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. Mauve: Measuring the gap between neural text and human text using divergence frontiers. *Advances in Neural Information Processing Systems*, 34:4816–4828, 2021.
- [58] Vitou Phy, Yang Zhao, and Akiko Aizawa. Deconstruct to reconstruct a configurable evaluation metric for open-domain dialogue systems. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4164–4178, 2020.
- [59] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- [60] Deven Santosh Shah, H Andrew Schwartz, and Dirk Hovy. Predictive biases in natural language processing models: A conceptual framework and overview. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5248–5264, 2020.

- [61] Albert Xu, Eshaan Pathak, Eric Wallace, Suchin Gururangan, Maarten Sap, and Dan Klein. Detoxifying language models risks marginalizing minority voices. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2390–2397, Online, June 2021. Association for Computational Linguistics.
- [62] Hila Gonen and Yoav Goldberg. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 609–614, 2019.
- [63] Dirk Hovy and Anders Søgaard. Tagging performance correlates with author age. In *Proceedings of the 53rd annual meeting of the Association for Computational Linguistics and the 7th international joint conference on natural language processing (volume 2: Short papers)*, pages 483–488, 2015.
- [64] Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. Prosocialdialog: A prosocial backbone for conversational agents. In *EMNLP*, 2022.
- [65] Chiori Hori and Takaaki Hori. End-to-end conversation modeling track in dstc6. *arXiv preprint arXiv:1706.07440*, 2017.
- [66] Michel Galley, Chris Brockett, Xiang Gao, Jianfeng Gao, and Bill Dolan. Grounded response generation task at dstc7. In *AAAI Dialog System Technology Challenges Workshop*, 2019.
- [67] Prakhar Gupta, Shikib Mehri, Tiancheng Zhao, Amy Pavel, Maxine Eskenazi, and Jeffrey P Bigham. Investigating evaluation of open-domain dialogue systems with human generated multiple references. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 379–391, 2019.
- [68] Erinc Merdivan, Deepika Singh, Sten Hanke, Johannes Kropf, Andreas Holzinger, and Matthieu Geist. Human annotated dialogues dataset for natural conversational agents. *Applied Sciences*, 10(3):762, 2020.
- [69] Shikib Mehri and Maxine Eskenazi. Ustr: An unsupervised and reference free evaluation metric for dialog generation. *arXiv preprint arXiv:2005.00456*, 2020.
- [70] Tianyu Zhao, Divesh Lala, and Tatsuya Kawahara. Designing precise and robust dialogue response evaluators. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 26–33, 2020.

## A TBD Human Judgement Evaluations

The TBD-Q1-2023 was evaluated at both the dialogue- and turn-level by crowd workers on Amazon Mechanical Turk. Each dialogue was evaluated across 4 dimensions: coherence, informativeness, likability, and overall. Coherence (or Understanding) is a 5 item Likert scale with 1 representing “The entire conversation is incomprehensible” and 5 representing “The dialogue is very coherent and all the information conveyed is consistent”. Informativeness is a 5 item Likert scale with 1 representing “There is barely any information content in the dialogue, such as generic utterances, perfunctory responses, and repetition. Often the utterances in the dialogue are short. Dialogues that receive a rating of 1 for understanding/coherence” and 5 representing “Most of the utterances in the dialogue are long sentences with high information content, and all the information is correct”. Likability (or Engagingness) is a 5 item Likert scale with 1 representing “The content of the conversation is unattractive, and I don’t know how to continue the conversation; dialogues receive a rating of 1 for understanding/coherence” and 5 representing “The conversation is extremely attractive and I am eager to continue it”. Overall is a 5 item Likert scale with 1 representing “The overall quality is very low, the conversation is not fluent and there is no information” and 5 representing “The overall quality is excellent, the conversation is very smooth, the amount of information content is very high with great engagingness, it’s a very good response”.

The turn-level was evaluated for Grammatical Correctness (“The quality of the English grammar”), Appropriateness (“The response is appropriate given the preceding turn (Note: The appropriateness of a response is very subjective”), Content richness (“The response is informative, containing long sentences that include various entities (such as names of people, names of places or times), conceptual words (sky, dust, sorrow, etc.) or descriptive/emotional words (It hurts me, Lovely, etc.)”), and Relevance (“The response is related to the context of the dialogue and is good and smooth”). All items were on a 1 to 5 Likert scale, with 1 being lowest and 5 being highest (e.g., 1 = no grammatical correctness).

## B Additional TBD Dialogue-level Evaluations

Tables 3, 4, and 5 show the results of our human metrics predicting the Coherence, Informativeness, and Likability crowd sourced dialogue-level annotations on the TBD-Q1-2023 data set.

		Automatic Metric Alone	Agreeableness		Empathy		All Psych.	
			P	P+A	P	P+A	P	P+A
TBD	DialogRPT	.186	.115	.188	-.012	.182	.102	.190
	Mauve	.081	.115	.128	-.012	.068	.102	.124
	USL-H	.048	.115	.183	-.012	.073	.102	.173

Table 3: Dialogue-level results predicting the Coherence rating: Reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metrics as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables.

		Automatic Metric Alone	Agreeableness		Empathy		All Psych.	
			P	P+A	P	P+A	P	P+A
TBD	DialogRPT	.230	.112	.224	.061	.225	.121	.216
	Mauve	.141	.112	.166	.061	.142	.121	.159
	USL-H	-.005	.112	.116	.061	.110	.121	.161

Table 4: Dialogue-level results predicting the Informativeness rating: Reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metrics as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables.

## C Additional Data Sets

**DSTC6 (D6)** is dialogue data collected from Twitter users for customer service for 40,000 context-response pairs [65]. The dialogue context was evaluated using 10 Turkers on a 5 point Likert scale based on the relevance of the response.

**DSTC7 (D7)** is conversation data extracted from Reddit conversation threads [66]. The dataset contained 3 million conversational responses and 20 million facts. The dialogue context was evaluated by crowdsourced annotators using a 5 point Likert scale based on the relevance and interest of the response.

**English As a Second Language (ESL)** consists of 200 different three turn dialogue segments from an English learning site [1]. This dataset consists of 21 comparisons across 5 dialogue systems with a human baseline over 13K judgements.

**DailyDialog (GD)** is a dialogue dataset constructed using 100 dialogue contexts from the test set of the DailyDialog dataset [67]. The context-response pairs were annotated by Turkers using a 1 to 5 scale based on appropriateness.

**HUMOD (HU)** is a multi-turn movie dialogue dataset created from the Cornell Movie-Dialogs Corpus [68]. This dataset is human annotated on a 1 to 5 scale based on the relevance of human generated responses to the context of a fictional conversation on the movie script.

**Neural Conversation Model (NCM)** consists of 200 hand-crafted single turn prompts originally from the IT Helpdesk Troubleshooting dataset [1]. This dataset consists of 59 comparisons across 11 dialogue systems with over 33K pairwise comparisons.

**Persona-DSTC10 (PD10)** is an evaluation dataset for the DSTC10 challenge constructed from a sample of 500 dialogue segments from the PersonaChat dataset [1]. A total of 4,500 context-response pairs were rated using an automatic dialogue response evaluator.

**Topical-DTSC10 (TD10)** is an evaluation dataset for the DSTC10 challenge constructed from a sample of 500 dialogue segments from the TopicalChat dataset [1]. A total of 5,000 context-response pairs were rated using an automatic dialogue response evaluator.

**TopicalChat-USR (TP)** is a human evaluation dataset developed from the Topical-Chat dataset through the USR metric annotation [69]. The context-response pairs were annotated by Turkers using a different scales based on qualities of understanding (0-1), natural (1-3), maintains context (1-3), interesting (1-3), uses knowledge (0-1), and overall quality (1-5).

**PersonaChat-USR (UP)** is a human evaluation dataset developed from the PersonaChat dataset the same way as TopicalChat-USR [69]. The context-response pairs are annotated with the same USR annotation scheme as TopicalChat-USR using the same qualities and scales.

		Automatic Metric Alone	Agreeableness		Empathy		All Psych.	
			P	P+A	P	P+A	P	P+A
TBD	DialogRPT	.081	.043	.071	.054	.090	.063	.076
	Mauve	.021	.043	.035	.054	.046	.063	.047
	USL-H	-.014	.043	.033	.054	.077	.063	.081

Table 5: Dialogue-level results predicting the Likability rating: Reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metrics as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables.

**DailyDialog (ZD)** is a dialogue dataset constructed using 100 dialogue contexts from the test set of the DailyDialog dataset [70]. The context-response pairs were annotated by Turkers using a 5 point Likert scale based on appropriateness, language usage, relevance, and context.

**PersonaChat (ZP)** is a dialogue dataset consisting of context-response pairs collected from the test set of the PersonaChat dataset [70]. The appropriateness quality of the response were annotated by Turkers for each context.

## D Automatic Metrics

**BARTScore** is a metric that evaluates generated text that uses a pre-trained BART encoder-decoder model [52]. This formulates the generated text evaluation as a text generation problem by directly evaluating text by the probability of being generated from or generating other textual inputs and outputs.

**BERTScore** is an evaluation metric for text generation that compute the similarity of two sentences as a sum of the cosine similarities between pre-trained BERT contextual embeddings [53]. For dialogue systems, this computes the F1 scores by matching token embeddings in the human reference and system response.

**BLEURT** is text generation evaluation metric based on BERT that can model human judgements [54]. This uses pre-training scheme on BERT with synthetic data and fine-tune it to predict a human score with a mean squared error (MSE) loss when applied to dialogue systems.

**DialogRPT** is an ensemble model consisting of GPT-2 based models trained on human feedback data for tasks predicting the feedback and humanlikeness of responses [55]. This utilized a contrastive learning approach using confounding factors affecting feedback metrics of the dialogue.

**Prism** is a machine translation evaluation framework that uses a sequence-to-sequence paraphraser to score outputs conditioned on a human reference [56]. This uses a multilingual neural machine translation (NMT) model as a zero-shot paraphraser which was trained by treating the paraphrasing as a translation task.

**USL-H** is a dialogue evaluation metric that uses a composition of measurements for understandability, sensibleness, and likeability [58]. This uses models trained for valid utterance prediction (VUP) to determine validity along with next sentence prediction (NSP) and masked language modeling (MLM) models to measure sensibleness and likelihood of a response.

## E Human Judgements

Table 6 lists the human judgements used across the additional data sets used in the supplement. These crowd-sourced annotations are averaged across all evaluations for each turn or dialog, which depends on the data set (e.g., a single prompt may have multiple crowd-sourced evaluations for Appropriateness). The average evaluation is then used as the gold standard for each unit in the data set.

Judgement	Question Text	Likert Scale	Data Sets
Appropriateness	The response is appropriate given the preceding dialogue.	1-5	ESL, NCM, PD10, TD10, ZD, ZP
Relevance	The response content is related to the preceding dialogue.	1-5	EC, ED, EE, HU
Enjoy	How much did you enjoy talking to this user?	1-4	PC
Overall	What is your overall impression of the quality of this utterance?	1-5	D6, D7, GD, FC, FT, TP, UP

Table 6: Human evaluation metrics

## F Results

Tables 7 through 18 contain results for each data set. All tables report adjusted  $R^2$  from a linear regression model whose dependent variable is the human evaluation metric (described above). We create three models which contain varying sets of independent variables: (1) the automatic metric alone (“Automatic Metric Alone”), (2) the psychological metric alone (“P”), and (3) both the psychological and automatic metrics together (“P+A”). In all models, the independent variables are mean centered and standardized, so that the resulting standard deviation is equal to 1. Note that “All Psych.” contains all five psychological metrics: agreeableness, empathy, emotional entropy, emotion matching, and language style matching.

DSTC6 (D6; [65])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.080	.000	.080*	.001	.081***	.038	.103***	.002	.081	.006	.084***	.009	.110***
BERTScore	.195	.000	.196*	.001	.195	.038	.224***	.002	.196	.006	.200***	.009	.227***
BLEURT	.167	.000	.167	.001	.168***	.038	.183***	.002	.168	.006	.170***	.009	.187***
Prism ref.	.081	.000	.082	.001	.084***	.038	.093***	.002	.083	.006	.085***	.009	.101***
Prism Unref.	.024	.000	.024	.001	.026***	.038	.044***	.002	.026	.006	.029***	.009	.052***
Prism Context	.014	.000	.015	.001	.016***	.038	.042***	.002	.016	.006	.019***	.009	.050***

Table 7: DSTC6 data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

DSTC7 (D7; [66])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.087	.000	.088	.000	.087	.023	.095***	.001	.088	.000	.091***	.009	.099***
BERTScore	.130	.000	.131	.000	.130	.023	.140***	.001	.131	.000	.130	.009	.141***
BLEURT	.126	.000	.127	.000	.126	.023	.130**	.001	.127	.000	.126	.009	.131***
Prism ref.	.101	.000	.101	.000	.101	.023	.105**	.001	.101	.000	.101	.009	.106***
Prism Unref.	.021	.000	.021	.000	.021	.023	.028***	.001	.021	.000	.021	.009	.029***
Prism Context	.011	.000	.011	.000	.011	.023	.026***	.001	.011	.000	.012	.009	.028***

Table 8: DSTC7 data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

English As a Second Language (ESL; [1])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.182	.002	.182	.005	.182	.004	.182	.000	.185	.011	.192	.009	.191
BERTScore	.096	.002	.098	.005	.098	.004	.103	.000	.106	.011	.127	.009	.131*
BLEURT	.080	.002	.082*	.005	.084	.004	.081	.000	.081	.011	.094	.009	.098*
Prism ref.	.066	.002	.067	.005	.067	.004	.067	.000	.067	.011	.079	.009	.078
Prism Unref.	.011	.002	.012	.005	.015	.004	.010	.000	.011	.011	.020	.009	.023*
Prism Context	.007	.002	.009	.005	.011	.004	.013	.000	.013	.011	.036	.009	.040*

Table 9: ESL data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

DailyDialog (GD; [67])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.017	.005	.021	.000	.018	.005	.024*	.002	.020	.004	.025	.009	.038**
BERTScore	.116	.005	.119	.000	.116	.005	.115	.002	.117	.004	.119	.009	.122*
BLEURT	.121	.005	.123	.000	.121	.005	.135*	.002	.120	.004	.129	.009	.140**
Prism ref.	.014	.005	.018	.000	.015	.005	.022*	.002	.016	.004	.022	.009	.033**
Prism Unref.	.002	.005	.003	.000	.002	.005	.004*	.002	.001	.004	.002	.009	.012**
Prism Context	.018	.005	.025	.000	.018	.005	.017	.002	.024	.004	.018	.009	.027**

Table 10: DailyDialog (GD) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

HUMOD (HU; [68])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.078	.000	.078	.000	.078	.006	.086***	.013	.089***	.002	.078	.009	.099***
BERTScore	.120	.000	.121	.000	.120	.006	.121*	.013	.129***	.002	.122*	.009	.130***
BLEURT	.123	.000	.123	.000	.123	.006	.129***	.013	.132***	.002	.124	.009	.139***
Prism ref.	.062	.000	.062	.000	.062	.006	.075***	.013	.073***	.002	.062	.009	.088***
Prism Unref.	.006	.000	.006	.000	.006	.006	.008**	.013	.021***	.002	.009**	.009	.026***
Prism Context	.024	.000	.024	.000	.024	.006	.043***	.013	.033***	.002	.024	.009	.052***

Table 11: HUMOD (HU) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

Neural Conversation Model (NCM; [1])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.035	.001	.035	.000	.035	.009	.039*	.006	.041**	.001	.035	.009	.045
BERTScore	.013	.001	.014	.000	.013	.009	.026	.006	.019*	.001	.014	.009	.033
BLEURT	.019	.001	.019	.000	.020	.009	.032*	.006	.026**	.001	.019	.009	.039
Prism ref.	.019	.001	.020	.000	.020	.009	.030*	.006	.027**	.001	.019	.009	.038
Prism Unref.	.004	.001	.005	.000	.004	.009	.009	.006	.009**	.001	.005	.009	.015
Prism Context	.013	.001	.014	.000	.013	.009	.015*	.006	.015*	.001	.015	.009	.020

Table 12: Neural Conversation Model (NCM) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

Persona-DSTC10 (PD10; [1])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.060	.000	.060	.000	.060	.019	.071**	.012	.069**	.006	.062	.009	.078***
BERTScore	.025	.000	.025	.000	.025	.019	.041**	.012	.037	.006	.031	.009	.050**
BLEURT	.043	.000	.043	.000	.043	.019	.054***	.012	.053**	.006	.048	.009	.062***
Prism ref.	.019	.000	.019	.000	.019	.019	.035**	.012	.031	.006	.024	.009	.044**
Prism Unref.	.020	.000	.020	.000	.020	.019	.028**	.012	.030*	.006	.024	.009	.037**
Prism Context	.008	.000	.008	.000	.008	.019	.027**	.012	.018	.006	.012	.009	.033**

Table 13: Persona-DSTC10, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

Topical-DSTC10 (TD10; [1])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.063	.000	.063	.000	.063	.000	.063	.001	.063	.002	.063	.009	.062
BERTScore	.054	.000	.054	.000	.054	.000	.054	.001	.055	.002	.055	.009	.055
BLEURT	.049	.000	.049	.000	.049	.000	.049	.001	.049	.002	.049	.009	.049
Prism ref.	.036	.000	.036	.000	.036	.000	.036	.001	.037	.002	.037	.009	.037
Prism Unref.	.000	.000	.000	.000	.000	.000	.000	.001	.000	.002	.002	.009	.002
Prism Context	.005	.000	.005	.000	.005	.000	.005	.001	.005	.002	.006	.009	.005

Table 14: Topical-DSTC10, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

TopicalChat-USR (TP; [69])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.138	.017	.146	.000	.136	.110	.213*	.003	.136	.070	.168*	.009	.227**
BERTScore	.217	.017	.224	.000	.215	.110	.279*	.003	.215	.070	.247*	.009	.295**
BLEURT	.283	.017	.289	.000	.281	.110	.324*	.003	.282	.070	.298*	.009	.332**
Prism ref.	.177	.017	.186	.000	.175	.110	.236*	.003	.175	.070	.204*	.009	.252**
Prism Unref.	.204	.017	.209	.000	.202	.110	.238*	.003	.202	.070	.229	.009	.252
Prism Context	.018	.017	.033	.000	.017	.110	.117**	.003	.015	.070	.099**	.009	.162***

Table 15: TopicalChat data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

PersonaChat-USR (UP; [69])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.044	.003	.041	.003	.041	.130	.179**	.003	.042	.019	.060	.009	.183***
BERTScore	.109	.003	.107	.003	.106	.130	.230**	.003	.107	.019	.128	.009	.238***
BLEURT	.151	.003	.150	.003	.149	.130	.258**	.003	.148	.019	.161	.009	.260***
Prism ref.	.080	.003	.077	.003	.077	.130	.202**	.003	.077	.019	.097	.009	.208***
Prism Unref.	.098	.003	.095	.003	.097	.130	.159*	.003	.096	.019	.116*	.009	.168***
Prism Context	.023	.003	.020	.003	.020	.130	.141**	.003	.022	.019	.054*	.009	.157***

Table 16: PersonaChat-USR (UP) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$

DailyDialog (ZD; [70])													
	Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.	
		P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A
BARTScore	.148	.001	.149	.000	.147	.020	.192***	.004	.150	.013	.172**	.009	.204***
BERTScore	.171	.001	.171	.000	.170	.020	.196**	.004	.172	.013	.191**	.009	.208***
BLEURT	.208	.001	.207	.000	.207	.020	.245***	.004	.208	.013	.235**	.009	.258***
Prism ref.	.146	.001	.145	.000	.145	.020	.184***	.004	.146	.013	.169**	.009	.194***
Prism Unref.	.024	.001	.025	.000	.024	.020	.029**	.004	.029	.013	.031*	.009	.041***
Prism Context	.007	.001	.008	.000	.007	.020	.039***	.004	.010	.013	.024**	.009	.051***

Table 17: DailyDialog (ZD) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$



PersonaChat (ZP; [70])													
Automatic Metric Alone	Agreeableness		Empathy		Emo. Entropy		Emo. Matching		Lang. Style Matching		All Psych.		
	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	P	P+A	
BARTScore	.179	.002	.182	.001	.178	.000	.180	.025	.202***	.040	.200**	.009	.224***
BERTScore	.160	.002	.163	.001	.159	.000	.163	.025	.184***	.040	.187**	.009	.214***
BLEURT	.170	.002	.172	.001	.169	.000	.174	.025	.191**	.040	.191**	.009	.216***
Prism ref.	.132	.002	.134	.001	.131	.000	.134	.025	.157***	.040	.158**	.009	.184***
Prism Unref.	.002	.002	.004	.001	.001	.000	.001	.025	.027***	.040	.040***	.009	.060***
Prism Context	.026	.002	.027	.001	.025	.000	.025	.025	.042**	.040	.058***	.009	.072***

Table 18: PersonaChat-ZP (ZP) data set, reported linear regression adjusted  $R^2$  where  $P$  contains the psychological metric as the independent variable and  $P + A$  contains both the psychological and automatic metrics as independent variables. Benjamini-Hochberg corrected significance level: \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$